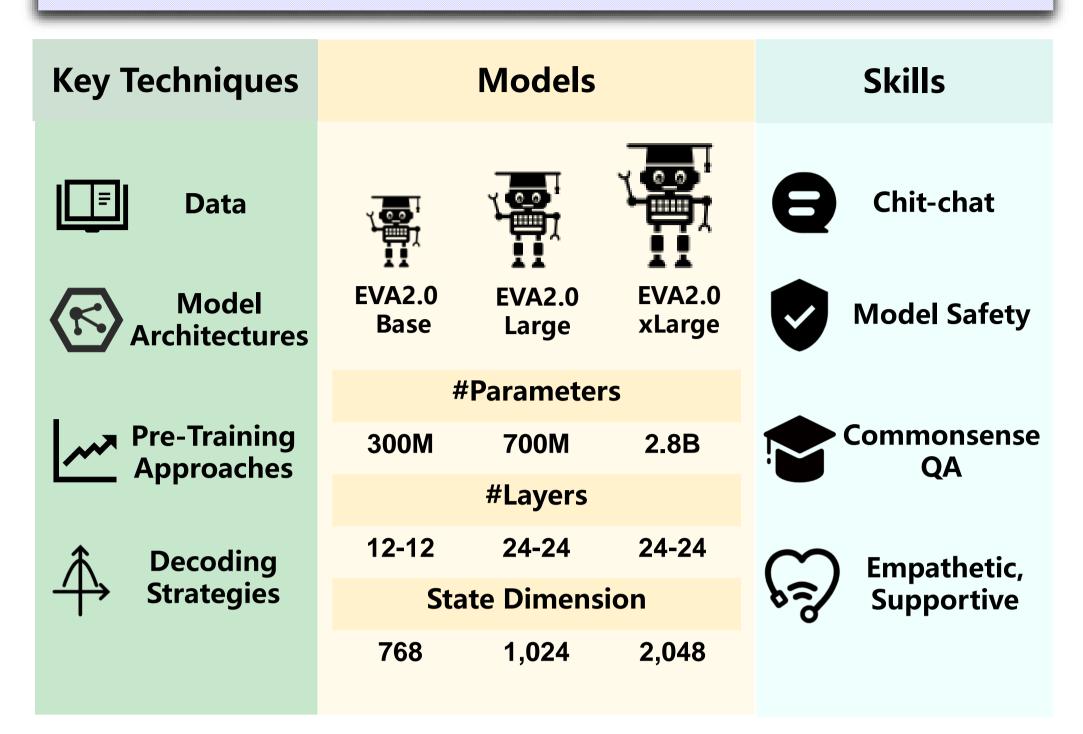
EVA2.0: Investigating Open-Domain Chinese Dialogue Systems with Large-Scale Pre-Training

Yuxian Gu*, Jiaxin Wen*, Hao Sun*, Yi Song, Pei Ke, Chujie Zheng, Zheng Zhang, Jianzhu Yao, Xiaoyan Zhu, Minlie Huang, Jie Tang The CoAl Group, Tsinghua University & BAAI

❖ Technical Framework



> Key techniques towards a human-like Chinese chatbot

- Data: Quality v.s. Scale
- Layer Number: Balanced v.s. Unbalanced Layers
- Role Information: w/ Role v.s. w/o Role
- Pre-training: Train from Scratch v.s. Further pre-training
- Decoding:
 - Base Method: Beam Search v.s. Sampling
 - Length Control: Length Penalty v.s. Min Length
 - Handling Repetition: w/ Control v.s. w/o Control

Data

Data Collection

• all data is publicly available

Weibo Comment	Weibo Repost	Zhihu
Baidu Zhidao	Subtitle	Story Dialog
Baidu Tieba	Douban	Public Dataset

> Data Quality Control

Coarse Filtering	Fine Filtering
context-level filtering sensitive word filtering	relevance score (BERT) fluency score (LM)
•••	•••

Data Analysis

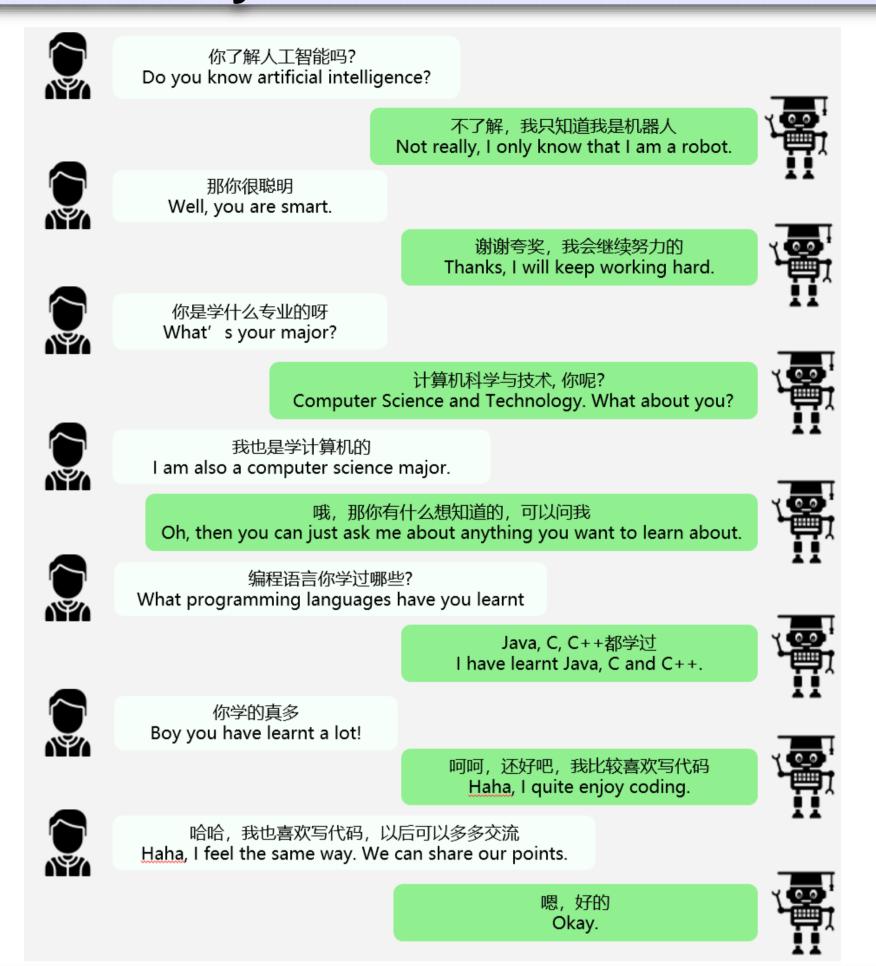
- 60GB high-quality dialogue pre-training dataset
- Basic statistics

Dataset	#Sess.	#Uttr.	#Token	Storage
WDC-Dialogue (Zhou et al., 2021)	1.4B	3.0B	78.3B	181GB
EVA2.0-dataset	0.4B	1.1B	22.4B	60GB

Quality statistics

Dataset	Relevance ↑	Fluency ↑	Entertainment ↓
WDC-Dialogue (Zhou et al., 2021)	55.2	-7,147	7.0%
EVA2.0-dataset	93.8	-3,237	6.2%

Case Study



Evaluation

> EVA2.0 significantly outperforms other open-source counterparts in both automatic and human evaluations.

> Strategies Comparison

Techniques	Model	F1	Single R-L	E-Turn B-4	D-4	F1	Multi R-L	-Turn B-4	D-4
Model	6-18 18-6 12-12 * +role	15.6 15.5 16.2 13.3	13.3 13.4 13.8 11.3	1.48 1.52 1.63 1.29	49.4 50.0 53.4 45.6	16.1 16.2 16.6 14.4	13.7 13.9 14.3 12.0	1.54 1.43 1.74 1.31	46.2 45.6 50.2 42.3
Pre-training	scratch * further	17.0 16.1	14.9 13.9	2.23 1.77	67.7 68.2	17.8 16.6	15.4 14.3	2.89 1.84	66.4 59.7
Decoding	greedy sampling beam search +sampling +len_penalty +no-repeat * +min_len	16.4 12.2 16.5 16.3 17.4 17.0 16.4	14.1 10.4 14.7 14.5 15.4 14.9 14.2	2.09 1.20 2.80 2.21 3.23 2.23 2.04	63.1 91.6 43.3 <u>75.4</u> 66.2 67.7 62.3	16.5 12.5 16.9 16.4 17.8 17.8	14.2 10.7 15.0 14.6 15.7 <u>15.4</u> 14.9	2.76 1.99 3.50 2.59 3.79 2.90 2.47	64.2 91.5 46.0 <u>73.2</u> 64.9 66.9 62.8

> Automatic Evaluation

Model	Single-Turn				Multi-Turn			
Model	F1	R-L	B-4	D-4	F1	R-L	B-4	D-4
CDial	9.9	8.6	0.67	61.2	11.9	10.3	0.88	63.9
EVA1.0	13.1	11.3	1.27	50.7	15.3	13.2	1.94	56.3
EVA2.0 _{Base}	16.2	13.8	1.63	53.4	16.6	14.3	1.70	50.2
$EVA2.0_{Large}$	16.4	14.0	1.67	55.8	17.0	14.9	2.03	53.9
$EVA2.0_{xLarge}$	17.0	14.9	2.23	67.7	17.8	15.4	2.90	66.9

> Human Evaluation

